

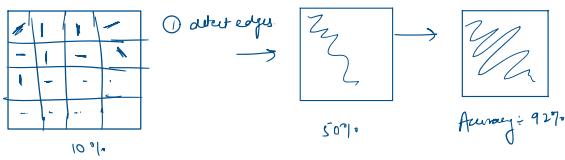
Convolutions Neural Networks

Computer Vision

- Object detection.
- Image classification.
- Neural style transfer.

 $64 \times 64 := (64 \times 64) \times 3$
 $\{rgb\}$

 $1024 \times 1024 := (1024 \times 1024) \times 3$
 $\approx 3 \text{ million.}$
(Very high)

Edge DetectionComputer Vision problems

- ① detect vertical lines: (|)
② detect horizontal lines: (—)

Vertical edge detection

* (Shift up 1 column after iteration k row after iteration j row)

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

Image: $(6 \times 6)_{x_1}$

Convolution: $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & -1 \\ 1 & 0 & 1 \end{bmatrix} = 4 \times 4 \text{ mat.}$
filter / kernel.

$(3 \times 1) + (0 \times 0) + (1 \times -1) + (1 \times 1) + (5 \times 0) + (8 \times -1) + (2 \times 1) + (7 \times 0) + (2 \times -1) = -5.$

$(0 \times 1) + (1 \times 0) + (2 \times -1) + (5 \times 1) + (8 \times 0) + (9 \times -1) + (7 \times 1) + (2 \times 0) + (5 \times -1) = -4$

* (Image - Vertical Edge detector)

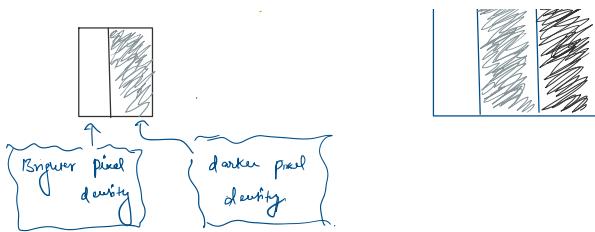
- tf. nn. Conv2d.

10	-	10	0	-	0
1	1	1	1	1	1
1	1	1	1	1	1
1	1	1	1	1	1
10	-	10	0	-	0

*

1	0	-1
1	0	-1
1	0	-1





$$\textcircled{1} \quad \text{pixel value} \propto \frac{\text{bright pixels}}{\text{dark pixels.}}$$

Horizontal Edge Detection

1	1	1
0	0	0
-1	-1	-1

More filters

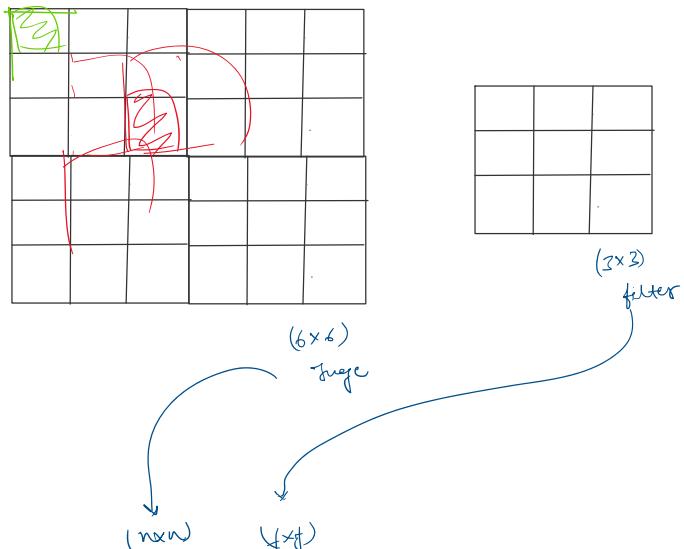
1	0	-1
2	0	-2
1	0	-1

3	0	-3
10	0	-10
3	0	-3

* Sobel filter * Scharr filter.

Vertical Edge detection.

Padding



$$(6,6) \xrightarrow[\text{f(3,3)}]{} (4,4) \xrightarrow{\text{Shrink}} (1,1)$$

$$(6,6) \xrightarrow[\text{f(3,3)}]{} (4,4) \xrightarrow{\text{Shrink}} (1,1)$$

- Two ways shrinking.
- Top left pixel [0,0] is used only once but centre pixel is overlapped many times.

- Top left pixel $[0,0]$ is used only once but centre pixel is overlapped many times.
- We are losing information from corner pixels.

Padding:

$$(6 \times 6) \iff (8 \times 8)$$

$(8 - 3 + 1) = 6 \times 6$ image (preserve image)

- padding of 0 :-
- padding = 1 :- $\frac{(n-f+1)}{n+2p-f+1}$

Valid and Same Convolutions

"Valid" :- No padding.

$$(n \times n) \rightarrow (f \times f) \rightarrow (n-f+1) \times (n-f+1)$$

"Same" :- Pad so that input size is same as output size.

$$(n+2p-f+1)$$

$$n+2p-f+1 = \cancel{x}$$

$$p = (f-1)/2$$

* f is usually odd.

$f \in 2n+1$

Strided Convolution

Stride \iff Step.

Stride = 2 (say).

Step key (Stride)

$$\left. \begin{array}{l} \text{padding: } p \\ \text{stride: } s \\ \text{image: } n \end{array} \right\} \left[\left[\frac{n+2p-f}{s} \right] + 1 \right] \times \left[\quad \right]$$

if this is not integer then floor ($\lfloor \cdot \rfloor$)

Cross-Correlation v/s Convolution

- ① Sometimes we flip the filter (vertically and horizontally) and this process is called "Convolution".
and process we perform is called cross-correlation.

$$f \rightarrow \begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 9 & 7 \end{bmatrix}$$

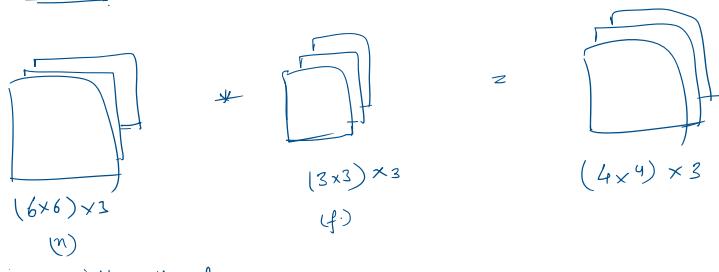
$$f^* \rightarrow \begin{bmatrix} 7 & 9 & -1 \\ 2 & 0 & 1 \\ 5 & 4 & 3 \end{bmatrix}$$

Convolution Over Volumes

$\downarrow \leftarrow \downarrow \rightarrow \curvearrowright$

Convolution Over Volumes.

① On RGB



height \times width \times channels



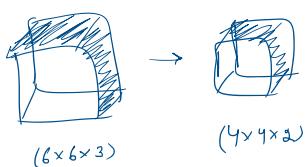
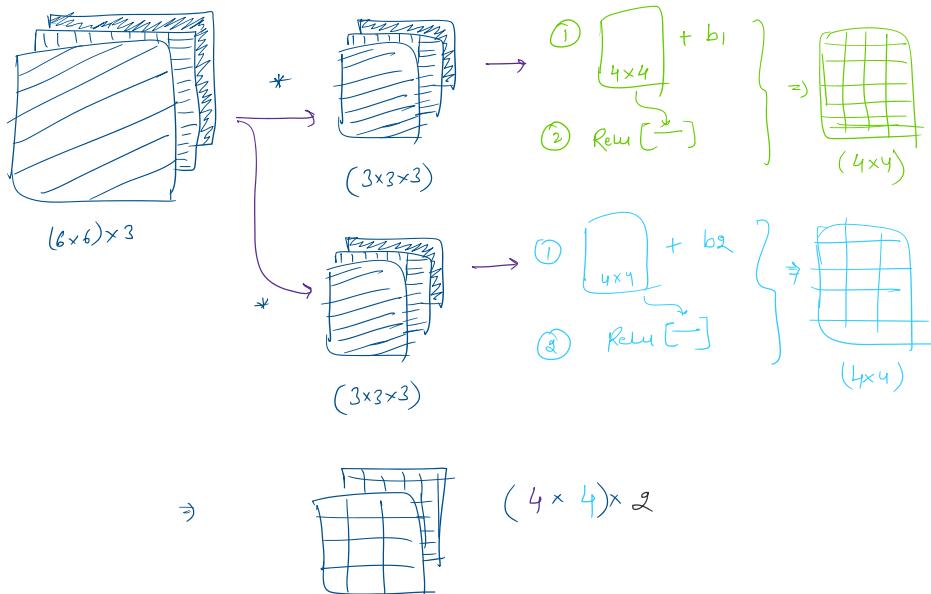
* Channels of image = Channels on filters.

$$(n \times n \times n_c) * (f \times f \times n_c) \Rightarrow (n-f+1) \times (n-f+1) \times n_c$$

$$(6 \times 6) \times 3 * (3 \times 3) \times 3. \quad 4 \times 4 \times 2$$

② Assuming stride=1, padding=None.

One Layer of CNN.



1 Layer of ConvNet.

- If layer l is a convolutional layer \Rightarrow

$f^{[l]}$ = filter size.

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = No. of filters.

Input : $(n_h^{[l-1]} \times n_w^{[l-1]} \times n_c^{[l-1]})$

Output : $(n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]})$

$$n^{[l]}_{(h/w)} = \left\lceil \frac{n^{[l-1]}_{(h/w)} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rceil$$

$f^{(c)}$ = filter size.
 $P^{(c)}$ = padding
 $S^{(c)}$ = stride
 $n_c^{(c)}$ = No. of filters.

$$\frac{\text{Size of filter}}{f^{(c)} \times f^{(c)} \times n_c^{(c)}}$$

Activations \Rightarrow

$$a^{(c)} \rightarrow n_{ht}^{(c)} \times n_w^{(c)} \times n_c^{(c)}$$

$$A^{(c)} \rightarrow m \times a^{(c)}$$

Weights \Rightarrow

$$(f^{(c)} \times f^{(c)} \times n_c^{(c)}) \times n_c^{(c)} \downarrow$$

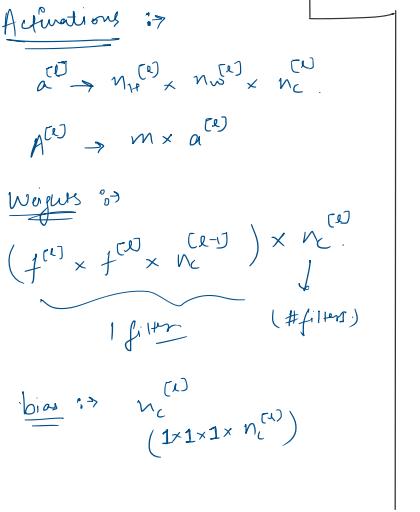
$\underbrace{f^{(c)} \times f^{(c)} \times n_c^{(c)}}_{\text{1 filter}}$ $\underbrace{n_c^{(c)}}_{(\#filters)}$

$$\text{bias} \Rightarrow n_c^{(c)} \quad (1 \times 1 \times 1 \times n_c^{(c)})$$

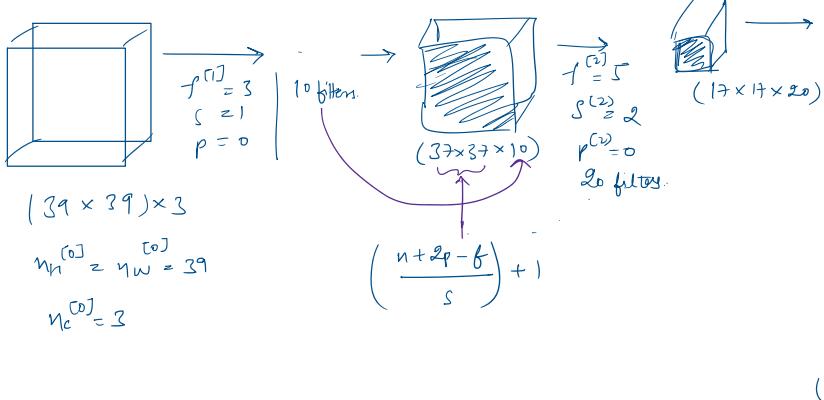
$$\text{Input: } (n_h^{(c)} \times n_w^{(c)} \times n_c^{(c)})$$

$$\text{O/p: } (n_h^{(c)} \times n_w^{(c)} \times n_c^{(c)})$$

$$n_h^{(c)} = \left\lceil \frac{n_h^{(c)} + 2P^{(c)} - f^{(c)}}{S^{(c)}} + 1 \right\rceil$$

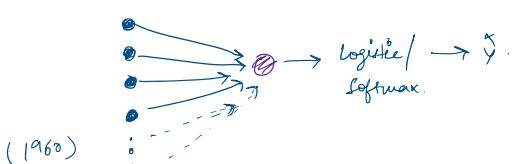


Simple Convolution N/w Example.



$$(7 \times 7 \times 40) \\ = 1960$$

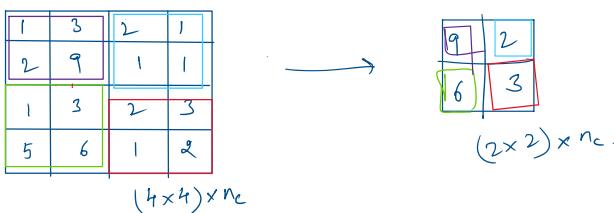
flatten



Types of N/w in CNN \Rightarrow

- Convolution. (CONV)
- Pooling Layer. (POOL)
- Fully Connected (FC)

Pooling layers : Max pooling



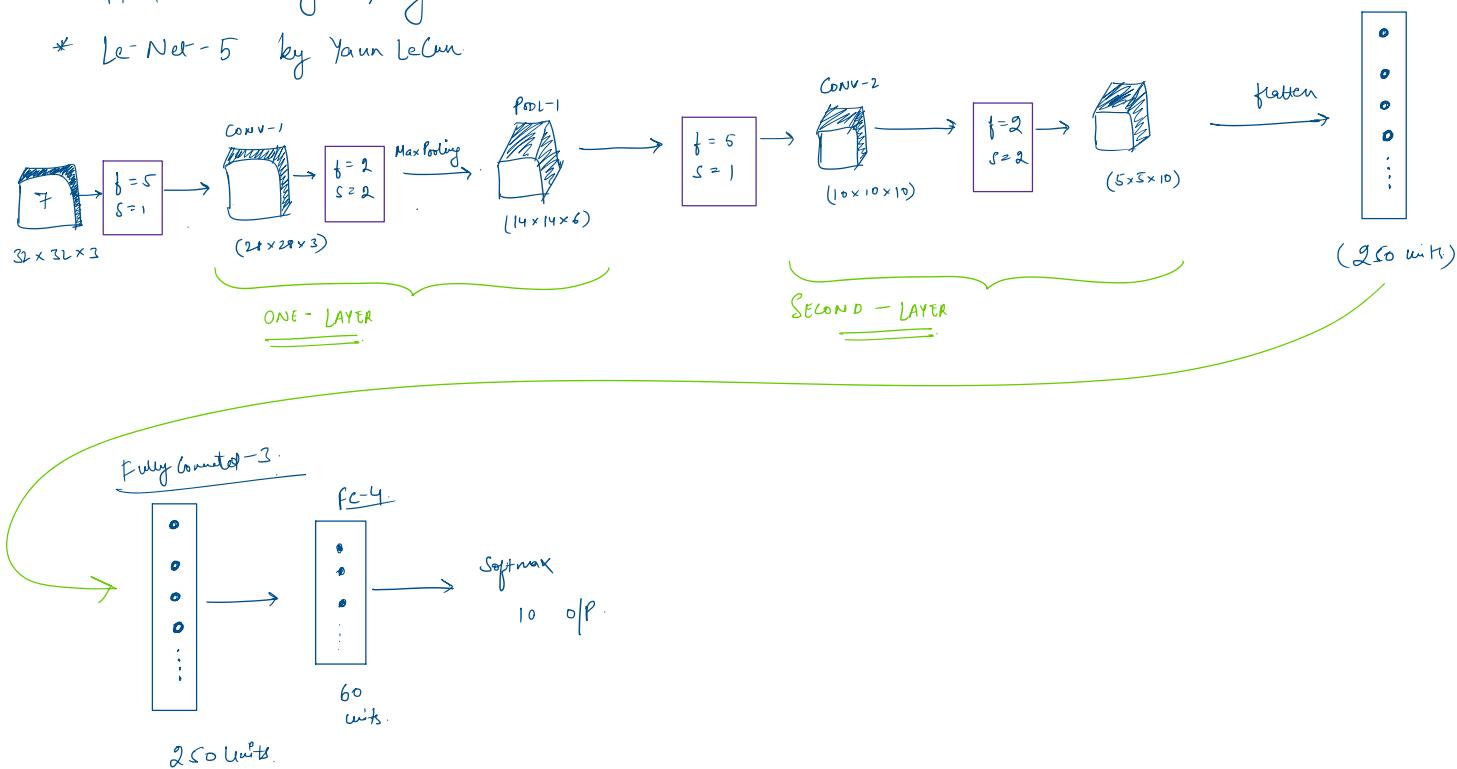
- Hyperparameters \Rightarrow
- * filter, $f = 2$
 - * stride, $s = 2$
 - * Take max over 2×2 regions.
 - * Takes no parameters.

$$\left(\frac{n+2p-f}{s} \right) + 1$$

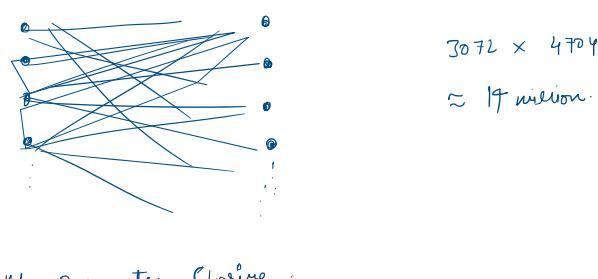
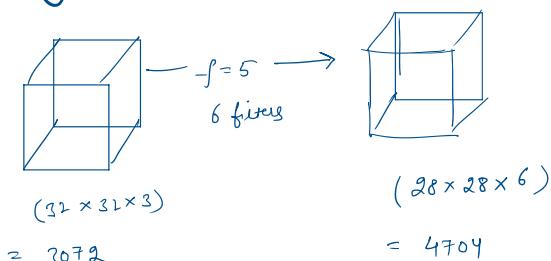
- Less used : Average pooling.

CNN Example

- Handwritten digit recognizer.
- * LeNet-5 by Yann LeCun.



Why Convolutions? ?



using the ReLU function.

#1 Parameter Sharing :-

A feature detector that's useful in one part of the image is probably useful in another part of image.

#2 Sparcity of Connections :-

In each layer, each o/p value depends only on a small number of units.

1 o/p unit is connected with $\frac{6}{36}$ units.
 \downarrow
 $(^n)$ $(n \times n)$

Classic N/w's.

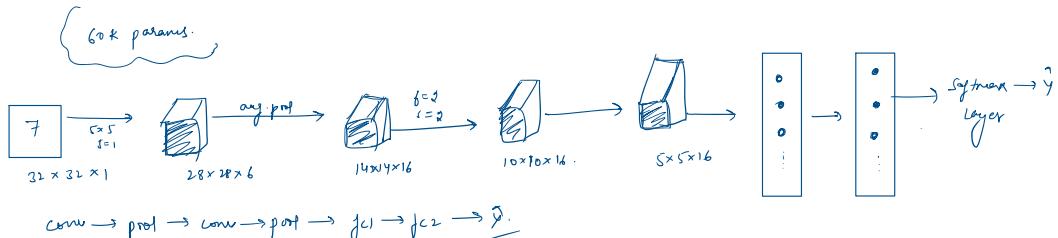
- LeNet-5
- AlexNet
- VGG-16

→ ResNet (Residual N/w).

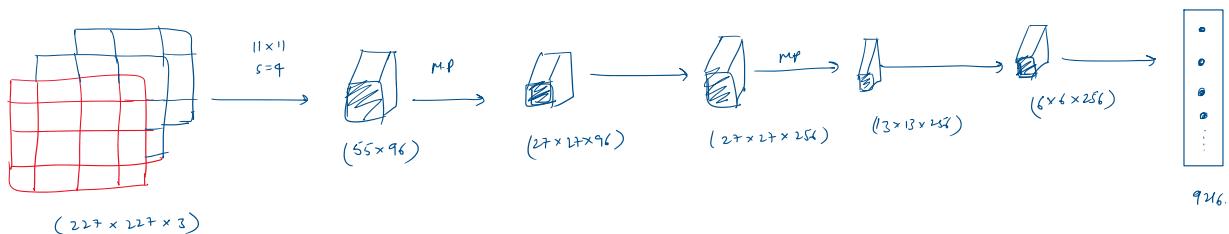
152 layer n.

→ Inception.

Le Net - 5



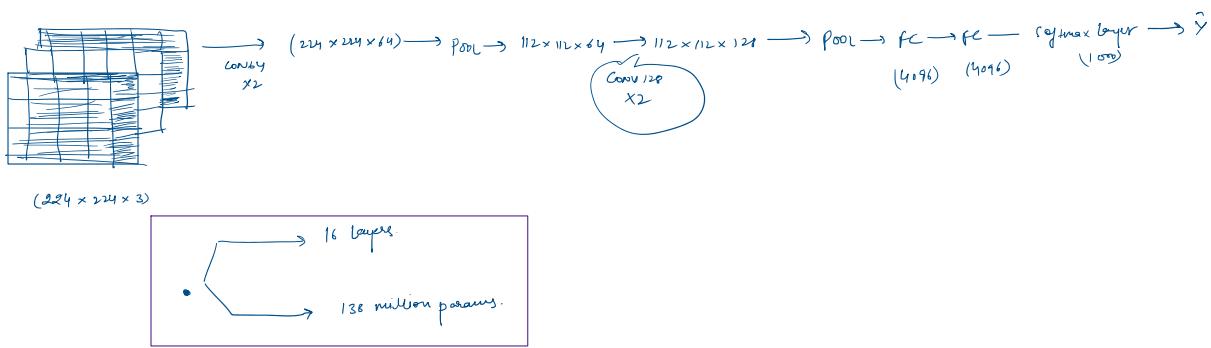
AlexNet



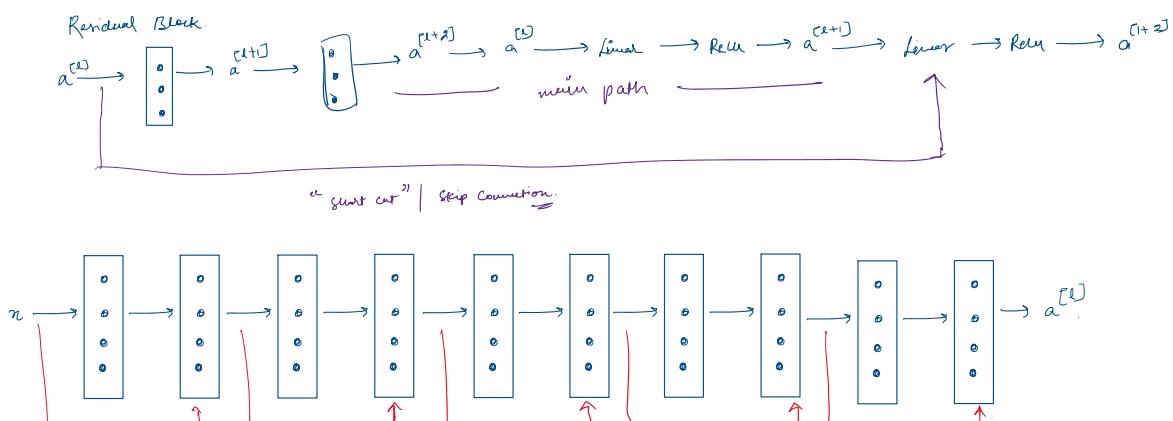
VGG-16

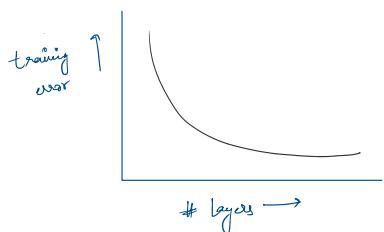
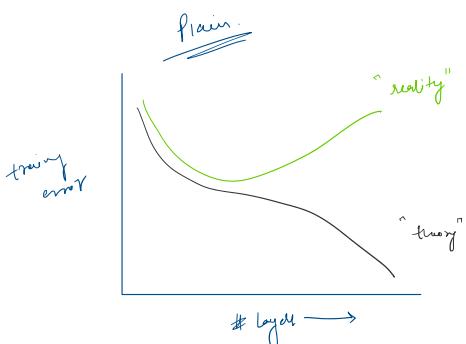
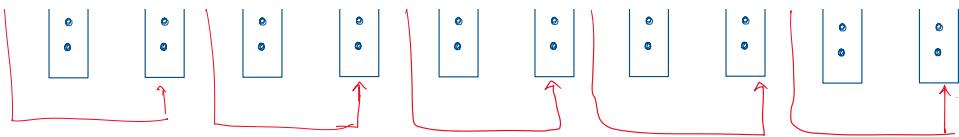
CONV = 3×3 filter, $s=1$, same.

MAX POOL = 2×2 , $s=2$.

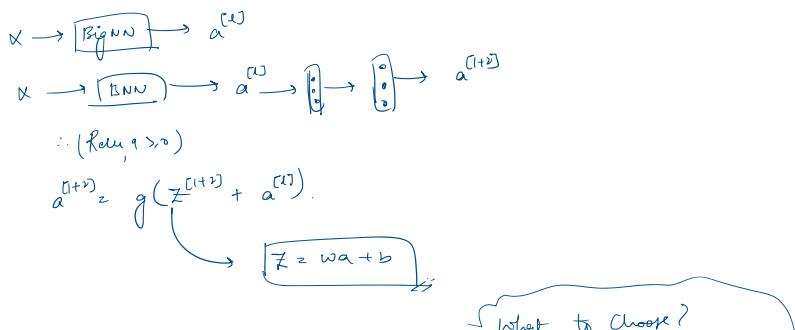


ResNets - Residual N/w's

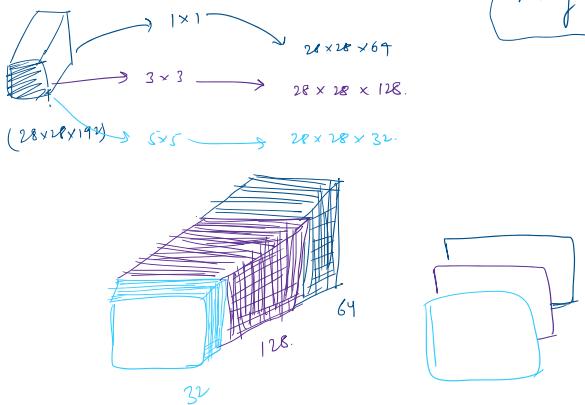




Why ResNets Work?



Inception N/W.



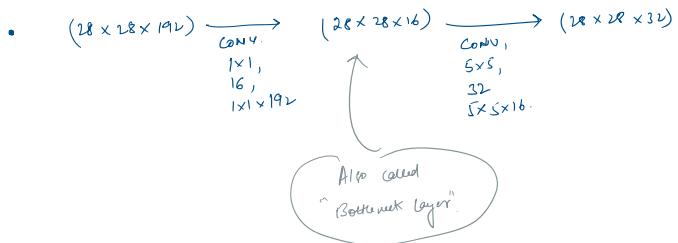
Computational Cost

$$(28 \times 28 \times 192) \xrightarrow{\text{CONV } 5 \times 5, \text{ same, } 32} (28 \times 28 \times 32)$$

$$(28 \times 28 \times 192) \times (5 \times 5 \times 192)$$

≈ 120 million.

Reduce cost by 1×1 convolutions!



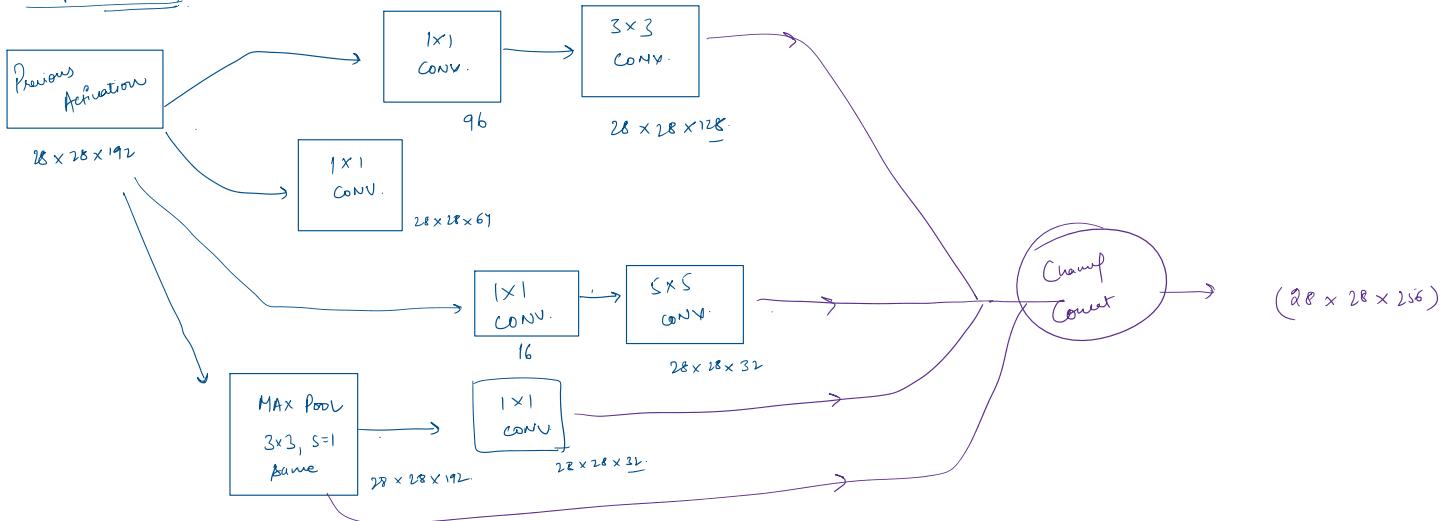
$$\text{Convolutional Cost.} \Rightarrow (28 \times 28 \times 16) \times 192 \approx 2.4 \text{ million}$$

+

$$(28 \times 28 \times 256) \times (5 \times 5 \times 16) \approx 10 \text{ million}$$

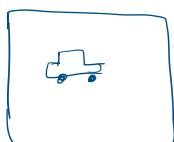
≈ 12.4 million

Inception Net



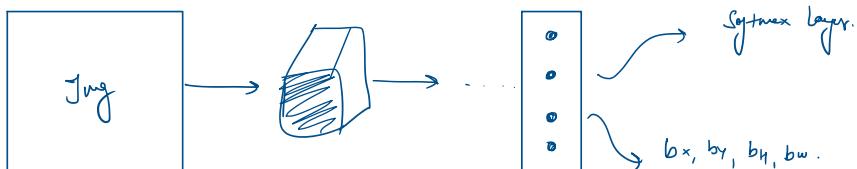
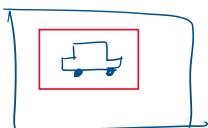
Object Localization

Image classification



1

Image classification with localization



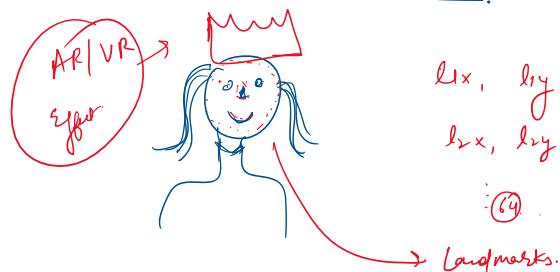
$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Probabilistic it is an object.

boundary boxes:

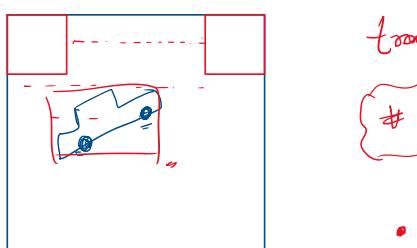
Class Names

Landmark Detection



Object Detection

Sliding window detection



Focus on with whole window.

increase the window size.

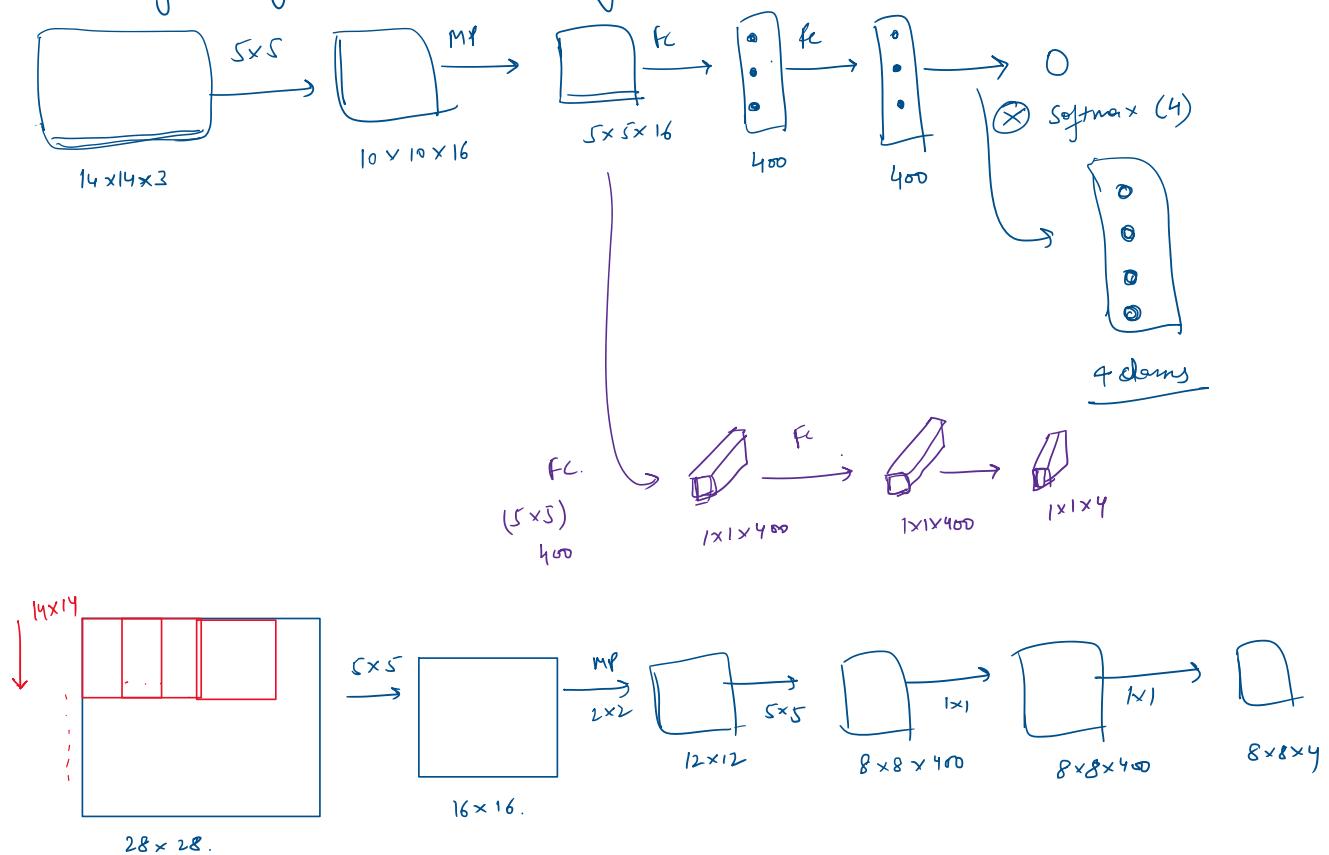
- Very Expensive Computational Cost



- Very Expensive Computational Cost

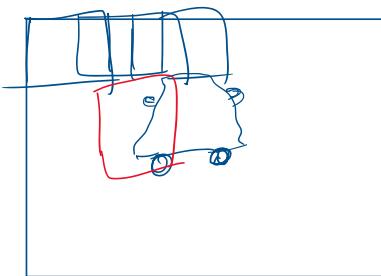
Convolution implementation of Sliding Window.

- Turning FC layers into Convolutional layers.



Bounding box prediction.

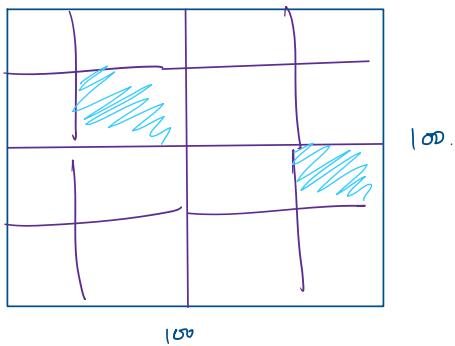
- Output accurate bounding boxes.



* None of the bounding box is Accurate!

YOLO Algorithm.

You only look once.



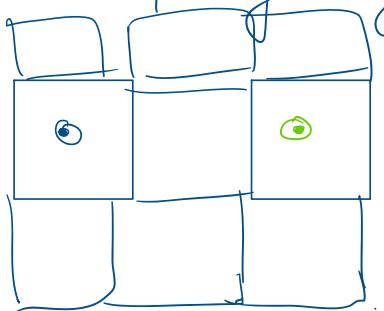
- * divide into grids.
- ~ Object to be detected.

Labels for training

for each grid cell -

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

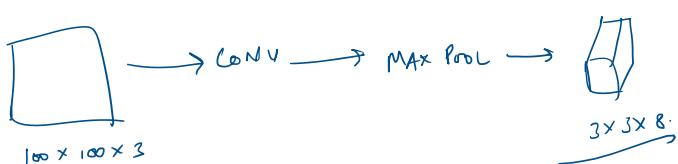
① Mid point of each grid cell is assigned if object is detected.



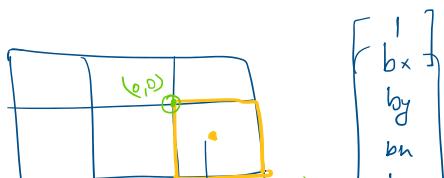
Total Volume of O/P \Rightarrow

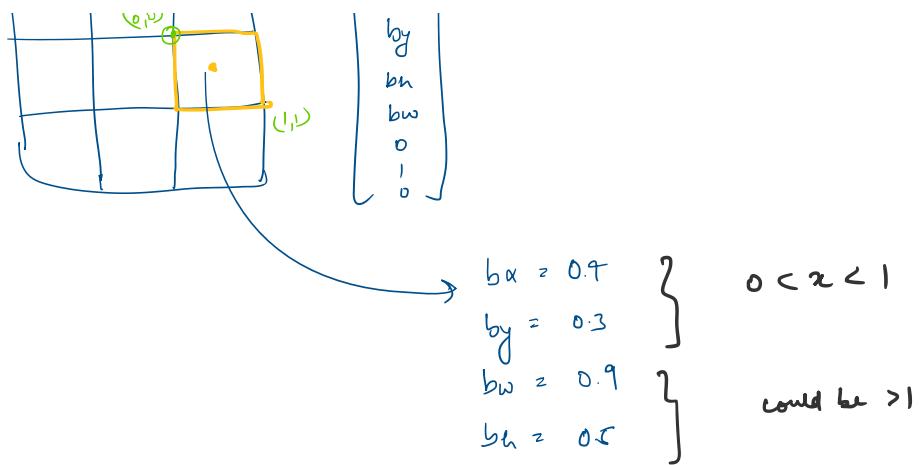
$$3 \times 3 \times 8$$

8 dimensional
o/p vector.



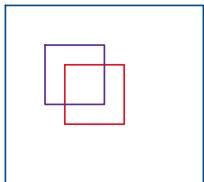
Specifying the bounding box.





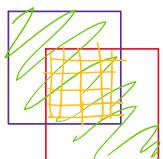
INTERSECTION OVER UNION.

Evaluating Object Localization



- ◻ \rightarrow ground truth bounding box.
- ◻ \rightarrow Algorithm Output.

IoU



Union

Intersection

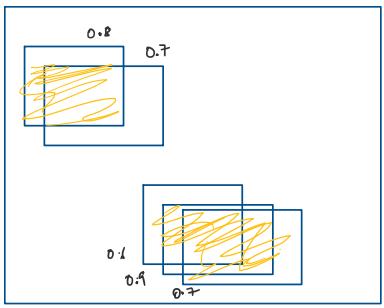
$$\frac{\text{Size of intersection} \#}{\text{Size of union} \#}$$

Correct if $IoU \geq 0.6$

- (*) IoU is the measure of Overlap.

NON-MAX SUPPRESSION.

- (*) Algorithm can detect same object multiple times, NMS makes sure that object is detected only once.

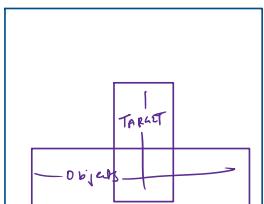


" $p(c)$ = Probability of Object."

- Discard all boxes with $p(c) \leq 0.6$
With the remaining boxes \Rightarrow
- Pick the boxes with maximum $p(c)$.
- Discard any remaining box with $IoU > 0.5$ with the box
Output in the previous step.

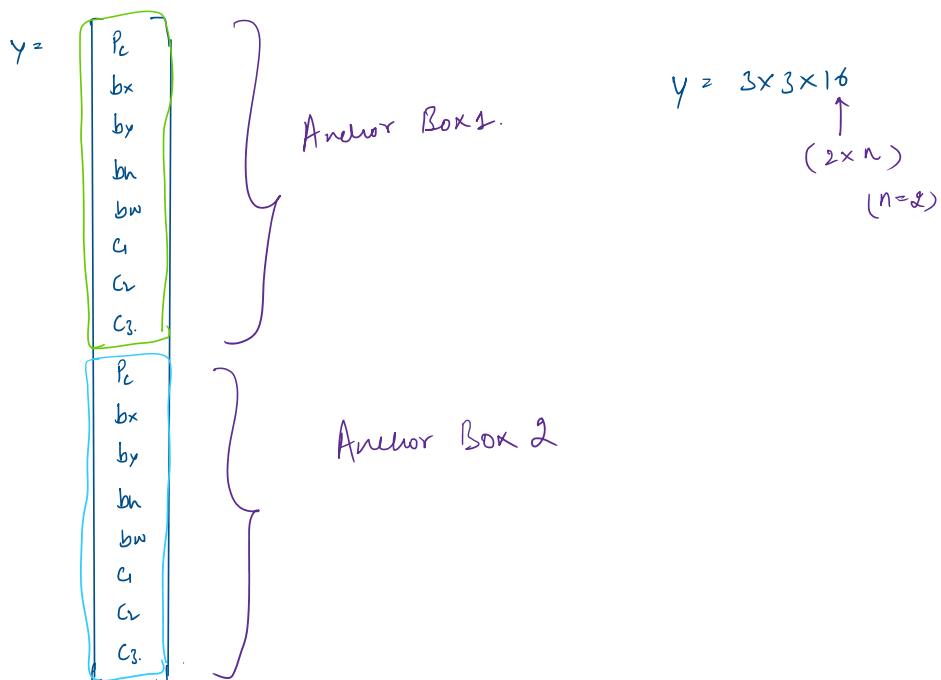
ANCHOR BOXES.

(X) Only one object in the gridcell can be identified



Predefine two (n) Anchor Boxes

- Anchor Box 1
- Anchor Box 2

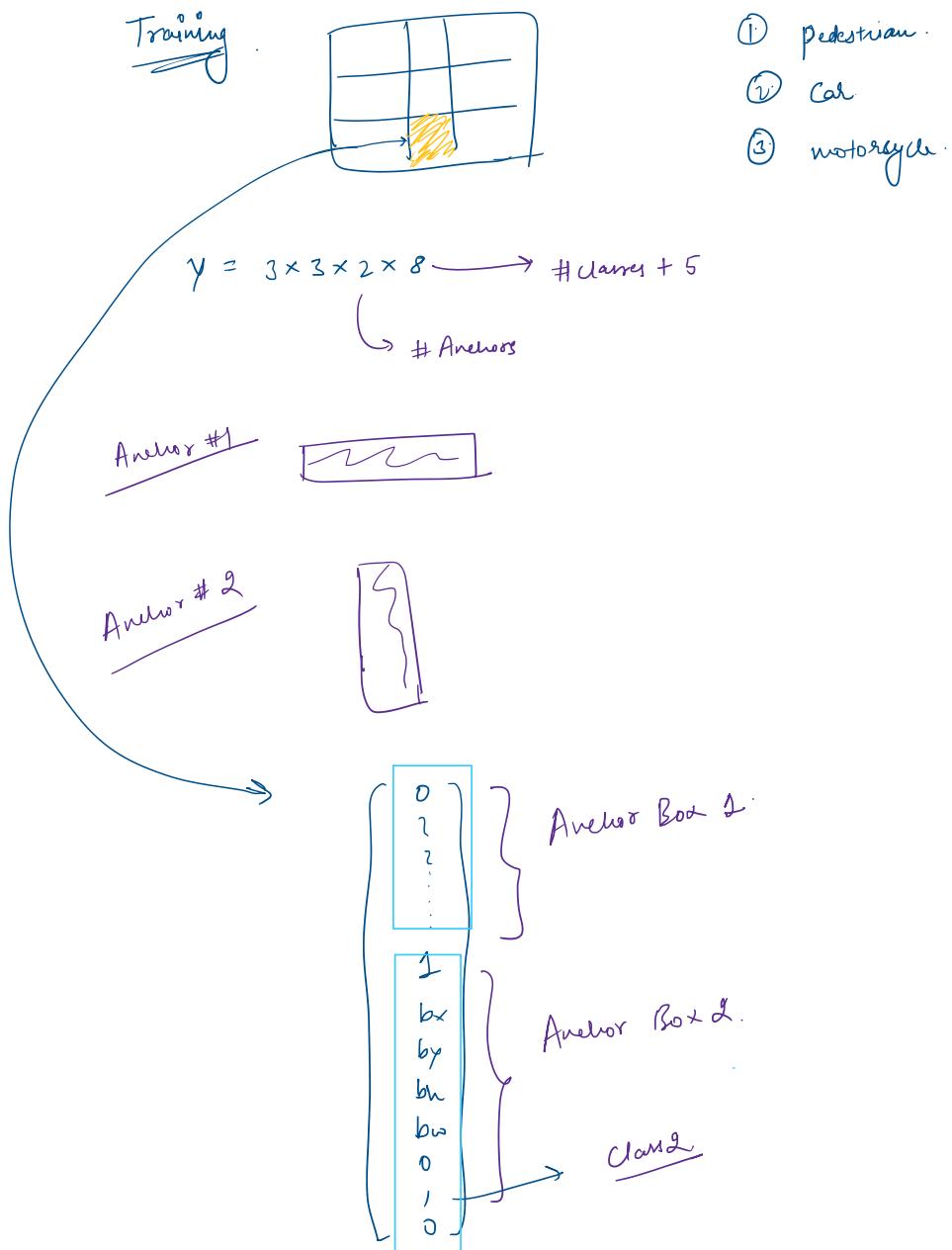


YOLO ALGORITHM : Putting it together.

Training



① Pedestrian



Outputting the non-max Suppnd O/P :-

- For each grid cell, get two predicted bounding boxes.
 - Get rid of low probability predictions.
 - For each class, use non-max suppression to generate final predictions.



FACE RECOGNITION.

① Verification

→ Input image.

→ O/p whether the image is of claimed person (0/1).

② Recognition.

→ Has database of K persons.

→ get an input image.

→ O/p Id of that person if in db.

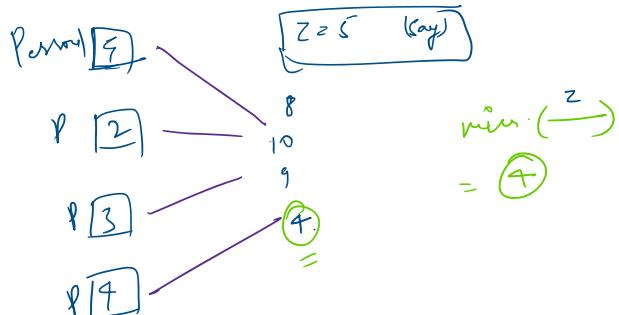
ONE SHOT LEARNING

"Learn by only single input image".

* Learn similarity fx.

$d(\text{img}_1, \text{img}_2) = \text{degree of difference b/w images img}_1 \text{ and img}_2$.

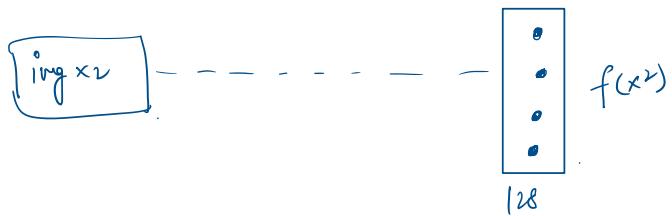
if $\begin{cases} d \leq z & \text{"Same"} \\ d > z & \text{"different"} \end{cases}$ $\quad \textcircled{z} = \text{Tau (some hyperparameter)}$



SIAMESE NETWORK.



Encoding of $\text{img}(x_1) = f(x_1)$.



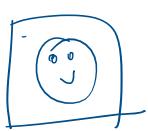
$$d(x_1, x_2) = \|f(x_1) - f(x_2)\|_2^2$$

=

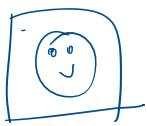
{ Comparing two images and calculating d is sometimes called
"Siamese Network".

"Deep Face"

Triplet Loss fn.

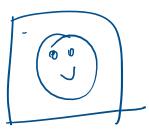


Anchor image.



Positive image.

$[z \ll]$



Anchor image.



Negative image.

$[z \gg]$

Triplet : { Anchor , A
+ve , P
-ve , N. }

- $\| f(A) - f(P) \|^2 \leq \| f(A) - f(N) \|^2$

$d(A, P)$ $d(A, N)$

$$d(A, P) - d(A, N) \leq 0$$

$$\| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 \leq 0 - \alpha$$

or

(hyperparameter)

$$\left[\| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 \right] + \alpha \leq 0$$

(Margin)

$\alpha = 0.2 \text{ (say)}$

Loss fx.

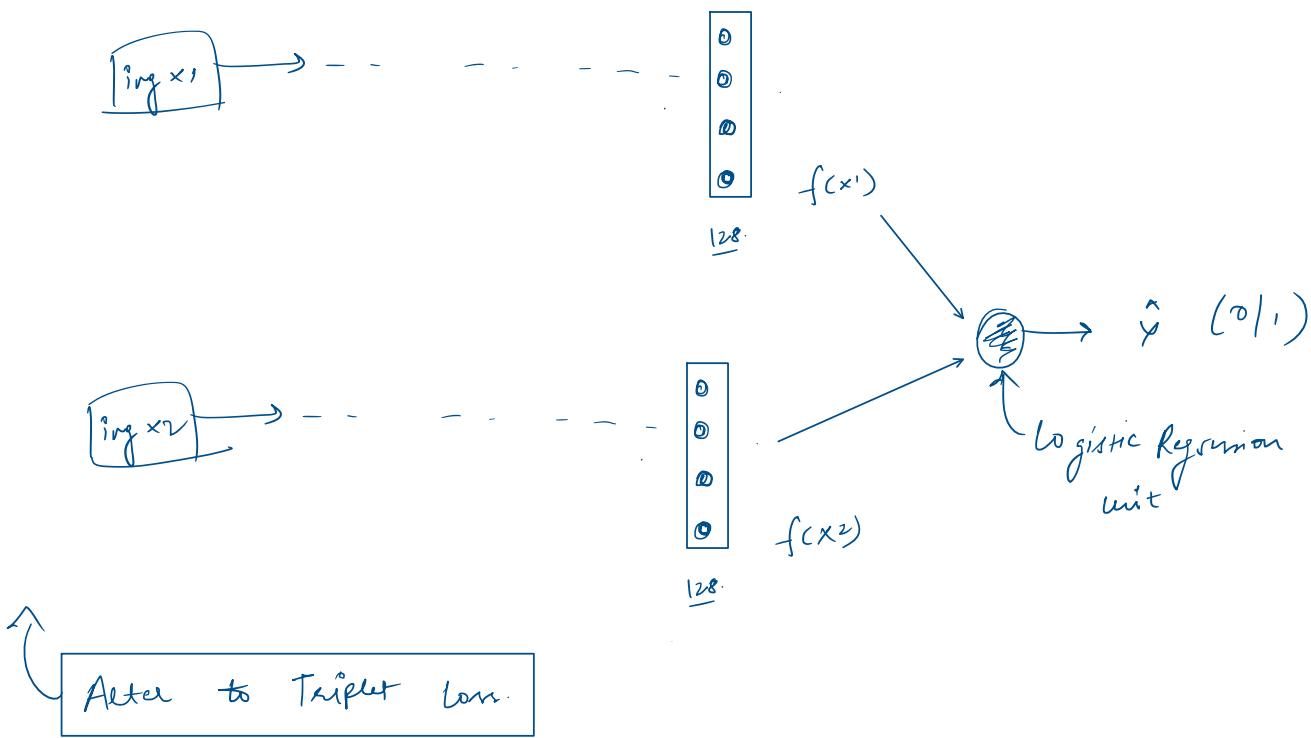
- given three images A, P, N .

$$\ell(A, P, N) = \max \left(\| f(A) - f(P) \|^2 - \| f(A) - f(N) \|^2 + \alpha, 0 \right)$$

≤ 0

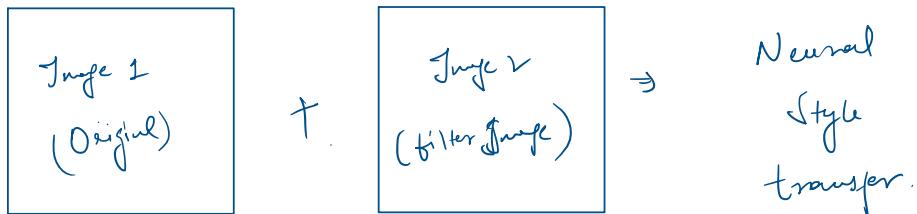
$$J = \sum_{i=1}^m \ell(A_i, P_i, N_i)$$

Face Verification & binary Classification.

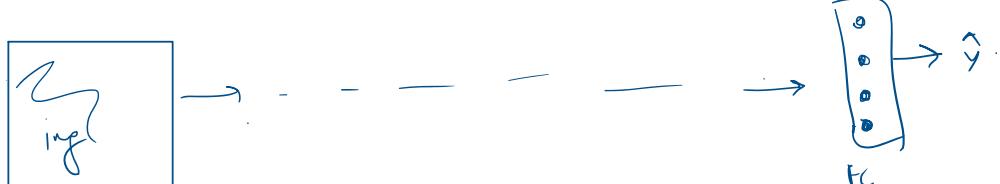


$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^i)_k - f(x^j)_k| + b \right)$$

NEURAL STYLE TRANSFER

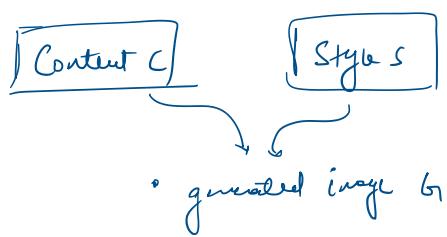


④ deep Convnets learning



Pick a unit in layer 1. Find the 9 image patches that maximize the units activation. Repeat for other units.

Cost fx.



$$J(G) = \alpha J_{\text{content}}(C, G)$$

↳ similarity.

+

$$\beta J_{\text{style}}(S, G)$$

Initiate G randomly,

$$G_1 : 100 \times 100 \times 3$$

$$G_t := G_{t-1} - \frac{\partial}{\partial G_t} J(G_t)$$

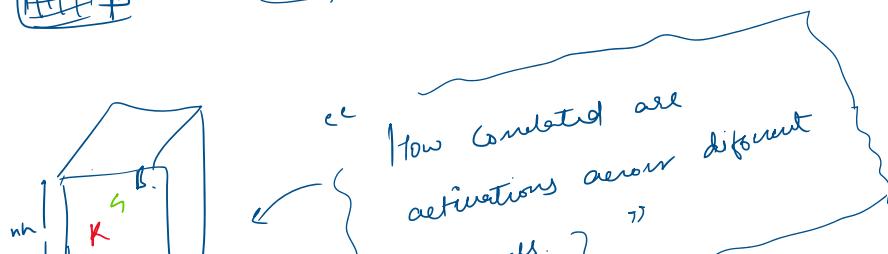
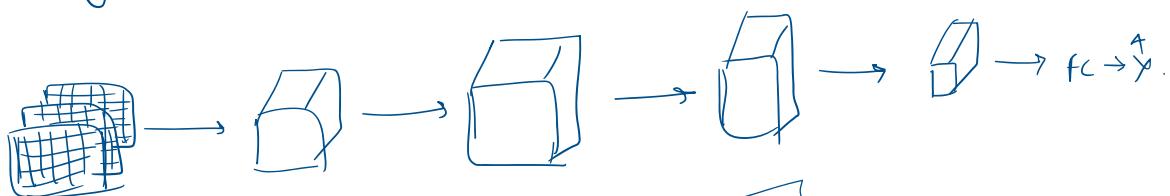
Content Cost fx.

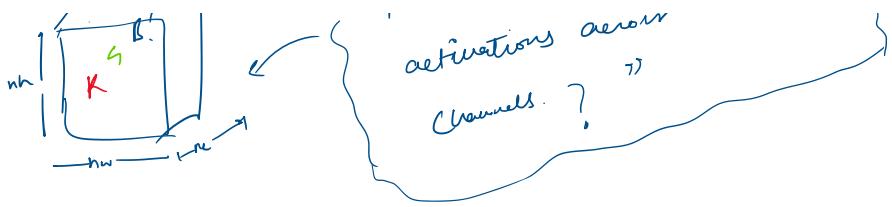
$$J(G) = \underbrace{\alpha J_{\text{content}}(C, G)} + \beta J_{\text{style}}(S, G)$$

- Say you use layer L .
- Use pre-trained ConvNet (say VGG).
- Let $a^{[L](C)}$ and $a^{[L](G)}$ be activation on layer L on the image.
- If $a^{[L](C)} \approx a^{[L](G)}$ are similar, both images have similar content.

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[L](C)} - a^{[L](G)} \|_2^2$$

Style Cost fx.





Style Matrix

Let $a_{ijk}^{(e)} = \text{activation at } (i, j, k)$

$G^{(e)}$ is $n_e \times n_c^{(e)}$.

$$G_{KK1}^{(L)(G)} = \sum_{i=1}^{n_H^{(L)}} \sum_{j=1}^{n_W^{(L)}} \left(a_{ijk}^{(L)(G)} \cdot a_{ijk}^{(L)(G)} \right).$$

(Correlation)

$$G_{KK1}^{(L)(S)} = \sum_{i=1}^{n_H^{(L)}} \sum_{j=1}^{n_W^{(L)}} \left(a_{ijk}^{(L)(S)} \cdot a_{ijk}^{(L)(S)} \right).$$

G = generated image.

S = style image

"gram matrix".

$$\text{JStyle}(S, G) = \| G^{(L)(S)} - G^{(L)(G)} \|_F^2$$

$$= \frac{1}{(2n_H^{(L)} n_W^{(L)} n_c^{(L)})^2} \sum_k \sum_{kl} (G_{kk1}^{(L)(S)} - G_{kk1}^{(L)(G)})$$

$$\text{J}(G) = \alpha \left[\text{JContent}(L, G) \right] + \beta \left[\text{JStyle}(S, G) \right].$$

1d and 3d generalization

$$\begin{array}{ccc} \text{2d image} & * & \text{2d filter} \\ (14 \times 14) & & (5 \times 5) \end{array} = (10 \times 10)$$

$$\begin{array}{l} \text{2d image} \\ (14 \times 14) \end{array} * \begin{array}{l} \text{2d filter} \\ (5 \times 5) \end{array} = (10 \times 10)$$

Same can be applied to 1d image.

ex. Ecg

electro cardio graphy

